

Optimization of the Center-of-Mass Trajectory with Reinforcement Learning

Alexandros Paraschos & Thibaut Munzer

Problem Statement

- Transfer the CoM from left to right foot
- Analytical solution is hard to find
- Working Approximation
- Room for improvement

-> RL!

Trajectory representation

- Dynamic Movement Primitives:

$$\gamma^2 \ddot{\mathbf{q}}_t = \alpha_q ((\mathbf{g} - \mathbf{q}_t) - \beta_q \dot{\mathbf{q}}_t) + f_t$$

Spring-Damper System:

$\mathbf{g} \dots$ goal-attractor
 α_q, β_q spring and damping coefficients

Forcing function: $f_t = \phi(\mathbf{z}_t)^T \mathbf{w}$

$\phi \dots$ Basis-Functions
 $\mathbf{w} \dots$ Shape-Parameters

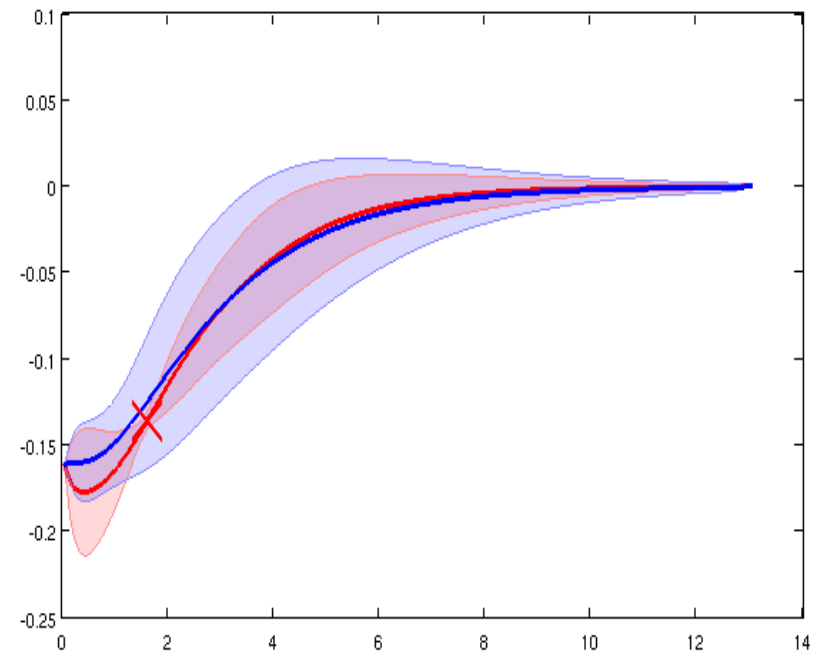
Phase variable: $\gamma \dot{z}_t = -\alpha z_t$

Replaces time, temporal scaling of the movement

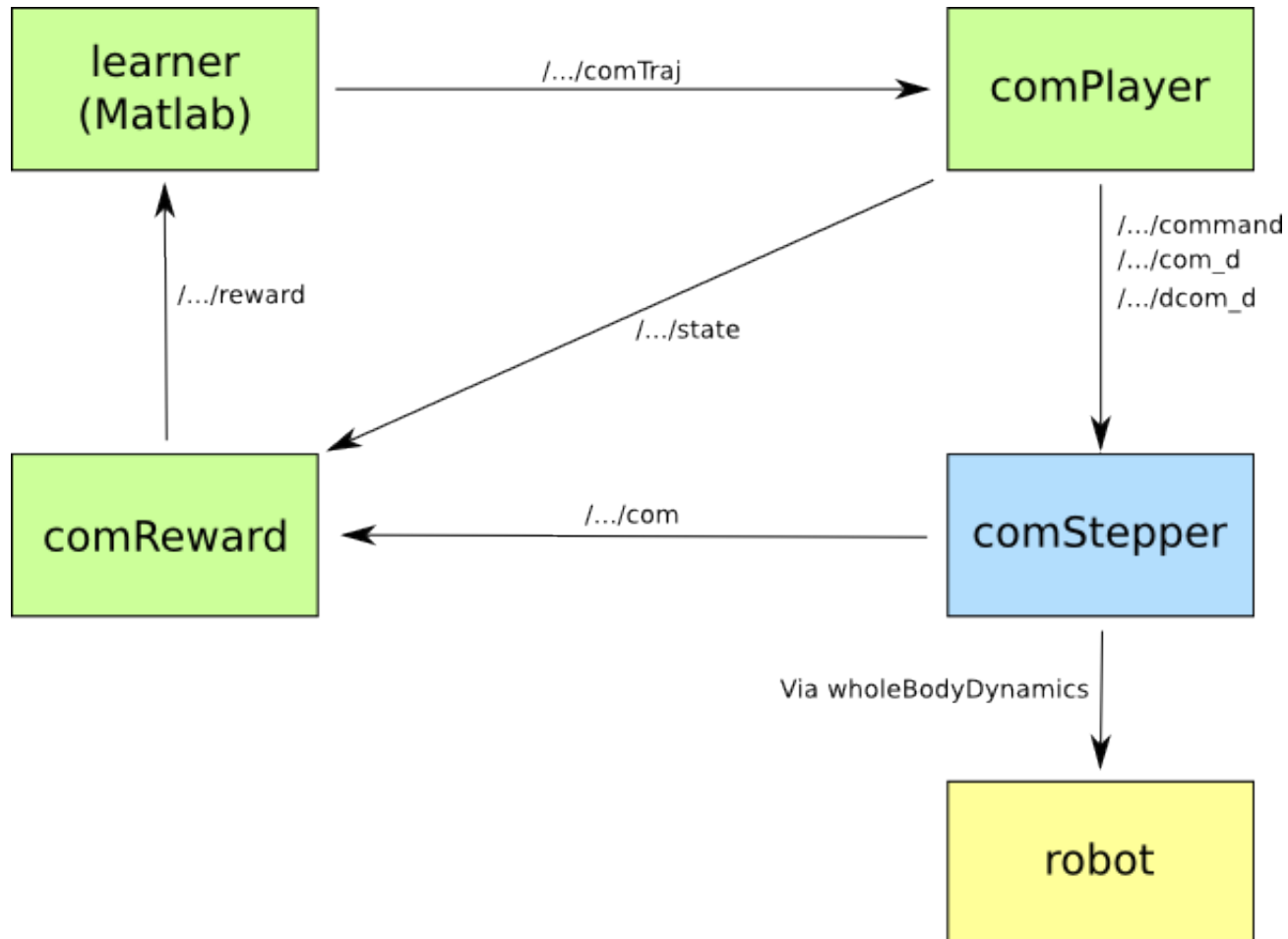
γ . temporal scaling factor

Reinforcement Learning Algorithm

- PoWER
 - Policy search
 - Maintain a distribution over the parameters
 - Samples new parameters
 - Weight each sample with the obtained reward
 - Update the distribution



RL Architecture



Experimental Setup

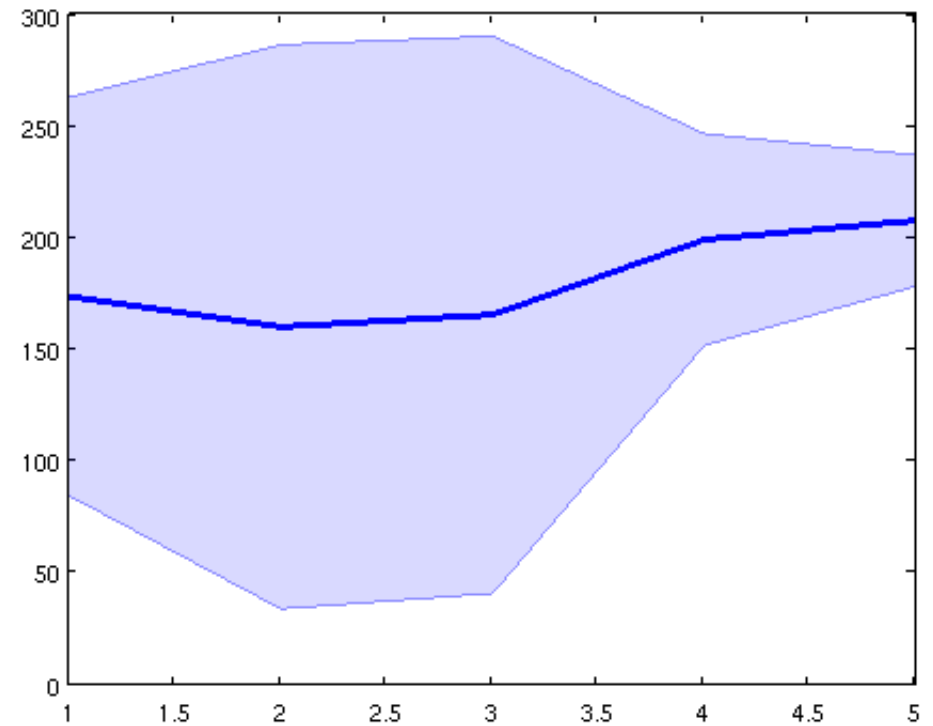
- Optimize the CoM projection (y , z)
 - 1 DMP for each dimension
 - 7 basis function per DMP
 - Plus the time coefficient
- PoWER algorithm
 - Episodic setup 10 samples per episode
- Initial policy: imitation of the analytical solution
- Reward
 - $R = 750 - \text{time not at objective} * 0.02$
 - Or $R = 0$ if the robot fall down during experiment

Results

- Not possible to run in simulation
- 5 Iterations on the real robot

	Reward
Demonstration	170
Demonstration (time scaled)	139
Best policy	231

Reward



Future work

- Evaluate for more iterations
- Average over multiple trials
- Learn to control the joint velocities
- Evaluate other RL algorithm
 - Bayesian optimisation, REPS, PI-BB
- Try different trajectory representations